# Tidy data and dates

Ronny A. Hernández Mora.

@RonnyHdezM
ronnyhdez
http://ronnyhdez.rbind.io/

# Workshop materials



https://github.com/ronnyhdez/curso_ciencia_datos_r

# What do we want from today's session?

- Understand what is tidy data
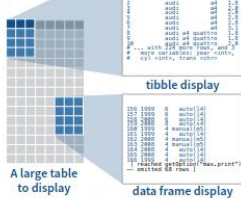- How to use tidyr
- Understand dates

# Tidy data

# RStudio cheatsheet

https://www.rstudio.com/resources/cheatsheets/

# What is tidy data?

| course | date | grade | estudents |
|--------|------|-------|-----------|
| matematica | 2020-05-28 | excelente | 34 |
| historia del arte | 2020-06-04 | regular | 20 |
| computacion | 2020-06-12 | bueno | 28 |

Each row is an **observation**

Each **column** is a **variable**

- Each **variable** must have its own **column**.
- Each **observation** must have its own **row**.
- Each **value** must have its own **cell**.

# Is this tidy data?

| course | date | grade | students/answers |
|---|---|---|---|
| matematica | 2020-05-28 | excelente | 34/20 |
| historia del arte | 2020-06-04 | regular | 20/18 |
| computacion | 2020-06-12 | bueno | 28/12 |

# Is this tidy data?

| course | date | grade | type | total |
|---|---|---|---|---|
| matematica | 2020-05-28 | excelente | estudiantes | 34 |
| matematica | 2020-05-28 | excelente | respuestas | 20 |
| historia del arte | 2020-06-04 | regular | estudiantes | 20 |
| historia del arte | 2020-06-04 | regular | respuestas | 18 |
| computacion | 2020-06-12 | bueno | estudiantes | 28 |
| computacion | 2020-06-12 | bueno | respuestas | 12 |

# How to make my data tidy?

```
pivot_longer( arguments )



                      pivot_wider( arguments )
```

https://tidyr.tidyverse.org

https://www.garrickadenbuie.com/project/tidyexplain/

Break  12:00 a 12:10

# Dates with lubridate

# RStudio cheat sheets



Dates and times with lubridate :: **CHEAT SHEET**

# Recursos

https://es.r4ds.hadley.nz/

# ¡Gracias!

Ronny A. Hernández Mora.

@RonnyHdezM
ronnyhdez
http://ronnyhdez.rbind.io/

ronny.hernandezm@gmail.com